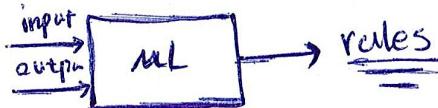


## machine learning:

- machine learning is branch of AI and computer science in which tries to learn and gain knowledge from input data
- to have great prediction  $\Rightarrow$  IBM definition ②
- without explicit coding



## Deep learning:

- IBM
- is subset of machine learning, which contains neural network
  - why Deep? using three or four or more layer
  - neural network is heartbreak
  - computing system  $\rightarrow$  tries to mimic human brain

## Computer Vision:

- Field of AI which gives ability to machines to process and extract meaningful information from Images/Videos like human.

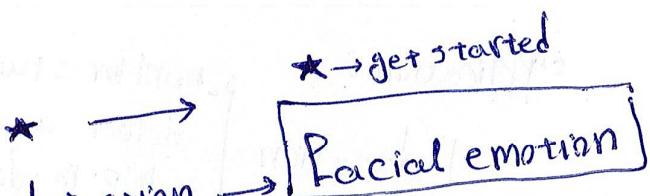
## Computer vision applications:

1) Image classification

2) object detection

3) Pattern recognition  $\rightarrow$  Face detection

4) Face detection



Facial emotion recognition  $\rightarrow$  FER (Face)

Explanation: FER is computer vision task aimed at identifying and categorizing emotional expression depicted on human face by analyzing Feature of Face: eyebrows, eye, mouth,

mapping them to

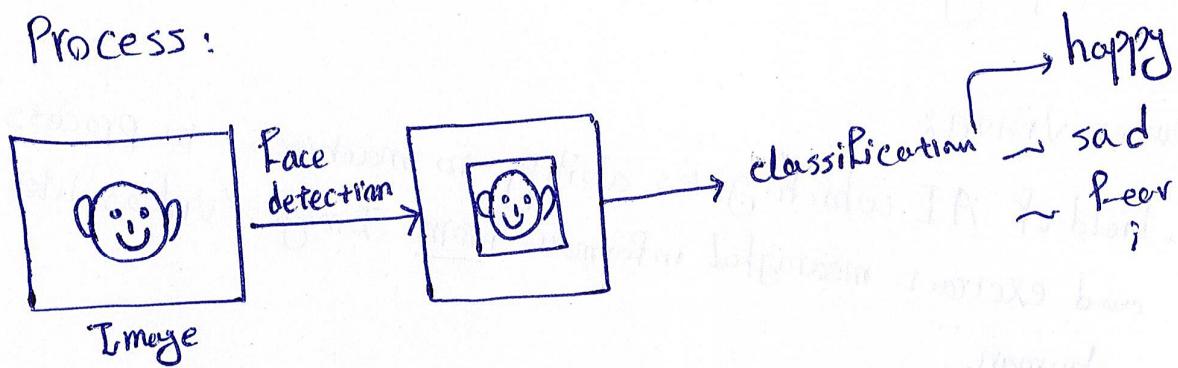
}	anger
}	Fear
}	surprise
}	sadness

- use Images / Videos

- real times

- Facial emotion is form of non-verbal communication

Process:



application:

Education

- {
  - monitor student's attention
  - detect a user's reaction to a topic
  - help to design an affective education system.

Health care

- {
  - detect Patient's Feeling.
  - observe patient condition.

Crime detection  $\Rightarrow$

{ spot shoplifter

## - Preprocessing

### 1) gray scale Image)

- they help to have less computation
- time saving
- without having any significant effect on accuracy specially on Facial emotion recognition.

### 2) resizing)

- Images are in different Dimension  $\Rightarrow$  they have to convert to the same size  $\rightarrow$  because we have introduced input layer.

### 3) Data augmentation }

- making more images based on feature which have been provided with current images.

### 4) Face Detection:

- The most important part of preprocessing
- algorithms/structures do not need background
- Faces are more important.

### 4,1) Viola-jones algorithm:

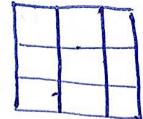
- humans can detect many objects easily (cars - trees, ...)
- this activity sounds crazy when it comes to computer.
- How does viola jones work?

★ Viola-Jones goal is to detect Face.

- to detect Face:

1) they look at the Frontal view of Face.

2) algorithm chunk image into section  $\Rightarrow$



3) examine each section and try to see if it sees any features of Face.

4) if it doesn't see any Features  $\rightarrow$  move to next section.

4.1) Features are  $\Rightarrow$  eyes, mouth, nose, cheeks, ...

★ Face can be considered as map and Face Features are landmark on maps.

★ algorithm is like a person trying to find the area that has all the landmarks in one location.

★ every time the algorithm detects a Face it marks where it found it on the pictures and keeps repeating. Process until it has checked whole images.

★★★ when the algorithm has finished scanning the Face:

- there are bunch of boxes that mark where Face located

- average on all location

## Viola Jones Face detection:

- Haar Features
- integral images
- adaboost
- cascading.

- it was trained with Faces and non-Face data to predict unseen

Face → real time - video - ..

- training Part | showing parts of image and telling the Face  
                  | is not Face

- goal → tries to train and gives ability to computer to understand what Face is and what is not Face by using picture data.

- once computer / program is trained → it can ~~be~~ extract Features  
Certain Feature

training process: - all certain Feature will be stored in a Rule

testing Part : 1) taking that important Rule

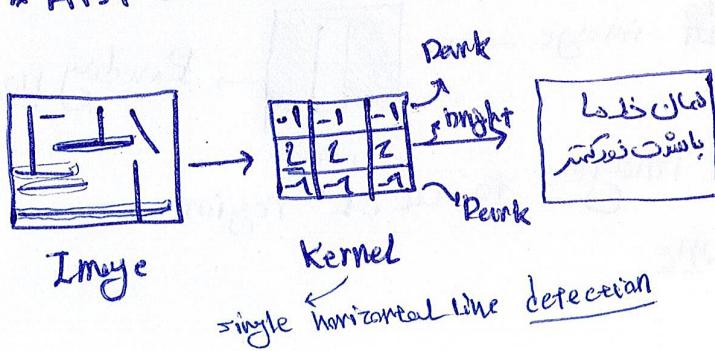
2) take new image <sup>in file</sup>

3) check all features <sup>in file</sup> and apply it on image

4) if it passes all the feature comparison → Face

- How are we going to extract Features?

- First of all lets talk about how to extract edge



\* by using this kernel:  
    Finding the single line (horizontal)

\* kernel moves on image

\* gives us the greater information  
    about horizontal line

## hair Features

- they are more similar to convolution kernels
  - ↳ try to find Feature and detect the presence of their Features.

- there 4 main Features?



Type 1



Type 2

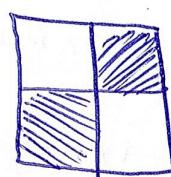


Type 3

Type 4



Type 5



- they are the same of kernel  $\rightarrow$  visualizing.

-   $\rightarrow$  black  $\rightarrow 1$   
 $\rightarrow$  white  $\rightarrow 0 - 1$

- when we apply this Feature to an image

- ↳ Function: ~~subtract~~ the pixel values ~~under~~ white region from pixels from black region

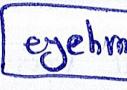
↳ output  $\rightarrow$  single value

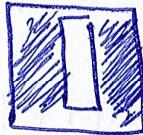
- consider another type of Features

- 2 - type 2  $\rightarrow$  apply it to image  $\rightarrow$  

$\rightarrow$  it will sum

all pixels ~~under~~ under black region and sum of all white pixels under white region  $\rightarrow$  subtract from each other

$\rightarrow$  single value comes out.  $\rightarrow$  like  eye brow

- 3 - type 3  $\rightarrow$  apply it ~~to~~ image  $\rightarrow$  

$\rightarrow$  finding nose

- we will apply these images ~~to~~ all of region.  
- finding the greatest one

~~Haar~~ →  
Haar-Features are able to find Facial landmark Feature (Nose, Mouth)  
by using those kernels.

- viola-jones uses  $24 \times 24$  window to move an image and apply  
Features.

→  → 2 pixel  $\Rightarrow$  it will apply to all pixels in image  $\rightarrow$  subwindow  $\rightarrow$  value extract  
If we increase size of kernel, what happens?

⇒ tip: increasing the size of each kernel  $\rightarrow$  Features 

Problem: 1) For every single window ( $24 \times 24$ )  $\rightarrow$  we have a huge amount  
of computation  $\rightarrow$  160,000 Features.   
    \*) all possible changes in  
    \*\*) position, scale, type  $\rightarrow$  all changes  
    \*\*) all changes in Haar-features

Idea: 1) eliminated all redundant Features.

    ↳ they are not useful

    ↳ other to do it? **adaboost technique**

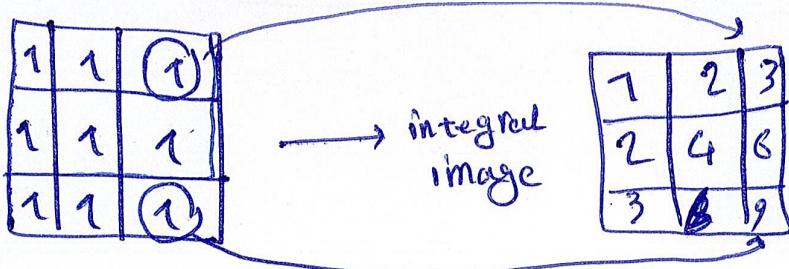
- before diving to adaboost  $\rightarrow$  there's important aspect.

Integral Images: \*) every single time, we need to sum up all white and  
black region  $\rightarrow$  Does not look efficient.  
- in an integral image the value at pixel  $(x,y)$  is sum of  
Pixels above and to left of  $(x,y)$

\*) every single part I need to sum up all the  
Pixels in black region and sum up white region

    ↳ takes time, computation.

    ↳ How to come up



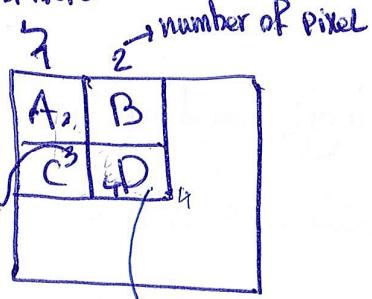
1	1	(1)
1	1	1
1	1	1

→ integral image

1	2	3
2	4	6
3	6	9

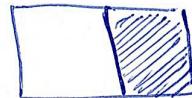
\* Integral image is ~~is~~ a new image matrix value.

\* when we are going to calculate the value of this patch



$$D = 1 + 4(2+3)$$

$$D = A + (A+B+C+D) - (A+C+A+B) \\ = D \Rightarrow$$



Imagine, we are going to calculate the value of D Patch in the integral image.  $\rightarrow$

1	2
3	4
5	6

Integral image allows for the calculation of sum of all pixels inside any given rectangle using only four values at corners of the rectangle.

- AdaBoost technique:

- ~~eliminate~~ eliminate the redundant feature.

- imagine we have calculated all computations

$\rightarrow$  160,000 feature

all of them are relevant? No.

relevant Feature

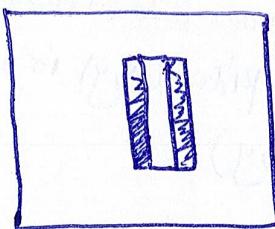


Image  
goal = finding nose

gives us great value

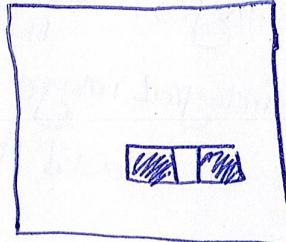


Image  
goal = finding nose

$\rightarrow$  we don't have nose here.

Irrelevant Feature

low. number

which one is relevant and which one is irrelevant's

2

- it will find certain number of features from all the 160,000 thousand Features

Laffer Lingay: It will give weight to these Features

and linear combination of all these Feature

is used to decide Face or no Face

$$F(x) = \alpha_1 \underbrace{f_1(x)}_{\substack{\alpha \neq 1 \\ \text{Strong classifier}}} + \alpha_2 \underbrace{f_2(x)}_{\substack{\alpha \text{ or } 1 \\ \text{Weak classifier}}} + \alpha_3 \underbrace{f_3(x)}_{\substack{\alpha \text{ or } 1 \\ \text{relevant Feature}}} + \dots$$

output  $\rightarrow 1 \rightarrow$  great  
 $\rightarrow 0 \rightarrow$  oh

\*) Cascading:

every 24\*24 window  $\rightarrow$  evaluate 2500 Features

after Performing adaboost: and then do the  
linear combination.

- instead of using 2500 Features  $\rightarrow$  we cascade

advantage: applying cascade on any given image  
we can detect Face or not-Face.

- Cascade is a group of stages.
- 2500 feature are grouped into several stages
- the job of each stage is determining whether a given sub window is definitely not a face or may be a face.

## MTcnn:

- multi-task cascade convolutional neural network
- MTcnn is model tool for face detection
- based on CNN → great performance - **robust**\*
- contains 3 important layers (levels)

- 1) Image is resized multiple times to detect face of different size
- 2) P-net → Proposal network scans → tries to find **First-detection** → it contains low threshold for detection  
→ so it detects so many false positive  
→ so it's slow
- 3) Proposed many regions (many false positive) → input for second network (R-net) → refine network  
→ tries to find bounding boxes → **by filtering detections**  
→ they are great \*
- 4) O-net → apply final refinement → to obtain great and accurate bounding box

## Comparison:

Features	Viola-Jones	MTcnn
speed	Very Fast	Fast
accuracy	Good	Very Good
Robustness	Bad	Very Good
using color	No	Yes
using GPU	—	Yes

Some advantages of a great detector:

- 1) Fast
- 2) accurate  $\rightarrow$  no False positive - no False negatives
- 3) Robust  $\rightarrow$  detects in different size, Position - rotation

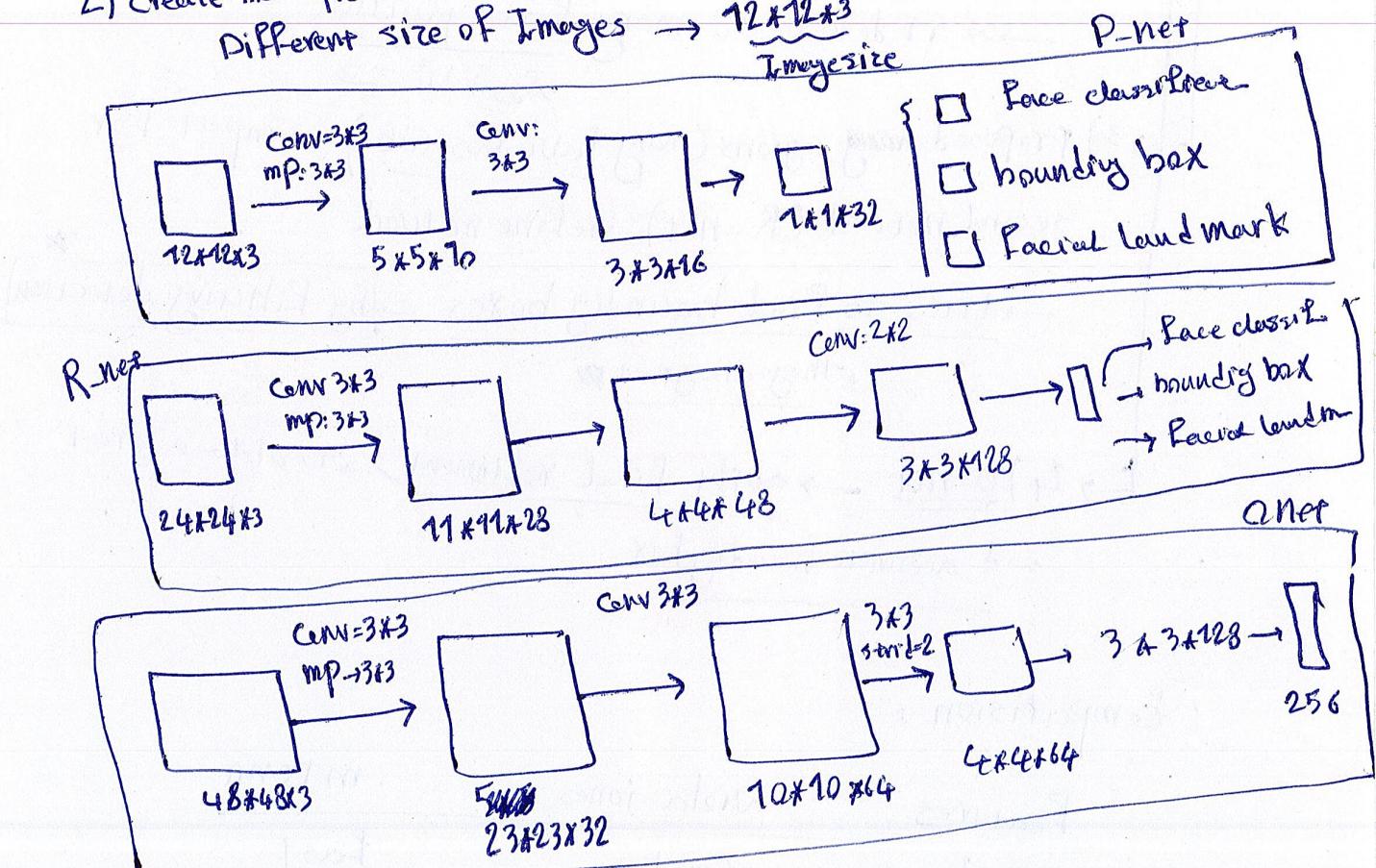
mTCNN is better than Viola-Jones

whole process:

1) Pass in image

2) create multiple scaled copies of the image  $\rightarrow$  Image Pyramid

Different size of Images  $\rightarrow$  Image size



mtchns:

↗

- Image Pyramid → making a copy of image in different sizes to detect different size of face.
- For each scaled copy → there is  $12 \times 12$  kernel that will go through every part of the image → aim, searching for face.  
From  $(0,0) \rightarrow (12,12)$
- this portion of image is passed → P-net
- P-net gives us the coordinate of bounding box if there is face
- this process will repeat  $(0+2a, 0+2b) \rightarrow (12+2a, 12+2b)$   
by shifting 2 pixel  $\Rightarrow$

Tip: small kernel in large image → detect small faces

- large kernel in small image → n large faces

- each network has its own parameter → P-net network parameter have been trained → so the output will be a accurate bounding box for each every  $12 \times 12$  kernel

- network gets a confidence score of each bounding box  
to make sure about face and delete those that are reliable or low confidence score.

- after collecting all boundary boxes with high confidence score  
we will have to standardize → convert to un-scaled image

- non-maximum suppression:

is method to reduce number of boundary boxes.

\* nms is conducted by first sorting boundary box by their confidence score.

- nms takes the largest boundary box instead of the one the network is most confident.

- calculating the area of each of the kernels,  
overlapping between each kernel and the kernel  
with highest score.

- delete that has high overlap. → make list of survived  
kernel.

stage 2) - converting all boundary boxes to 24x24

- normalize them

- R-net are more similar to PNet

- gives us near coordinate with high accuracy

- again, nms comes in to reduce number of boundary box

→ stages:

- before passing bounding boxes from R-net

→ they must convert to  $48 \times 48$   
similar to previous level → padding and then feed

→ these  $48 \times 48$  images go through O-net → CNN base

- outputs of O-net are slightly different from P-net, R-net.

- It produces three output

- 1) coordinate of bounding box → out[0]
- 2) coordinate of 5 facial landmark → out[1]
- 3) confidence score → out[2]

- again, we get rid of the boxes with lower confidence.

- standardize both bounding box and facial landmark coordinate  
to unscaled image

- another gate: NSM → only one bounding box for each face

## ★ short summary of the whole process ★

1) Pass in image

2) create multiple scaled copies of image.

3) Feed scaled image to P-net → CNN base (return bounding boxes)

4) Pick up P-net output

5) Delete bounding boxes with low confidence

6) Convert  $12 \times 12$  kernel coordinates to unscaled.

7) Non-maximum suppression → reduces the number of boundary box

8) Convert boundary boxes coordinates to unscaled.

9) reshape boundary boxes to square.

## Stage 2:

- 1) Pad-out bound boxes  $\rightarrow 24 \times 24$
- 2) Feed scaled image into R-net.
- 3) gather R-net output.
- 4) delete boundary boxes with low confidence.
- 5) NMS  $\rightarrow$  reducing number of boundary boxes.
- 6) convert them to unscaled image coordinate
- 7) reshape boundary boxes to square

## Stage 3:

- 1) Pad-out bound boxes  $\rightarrow \underline{48 \times 48}$

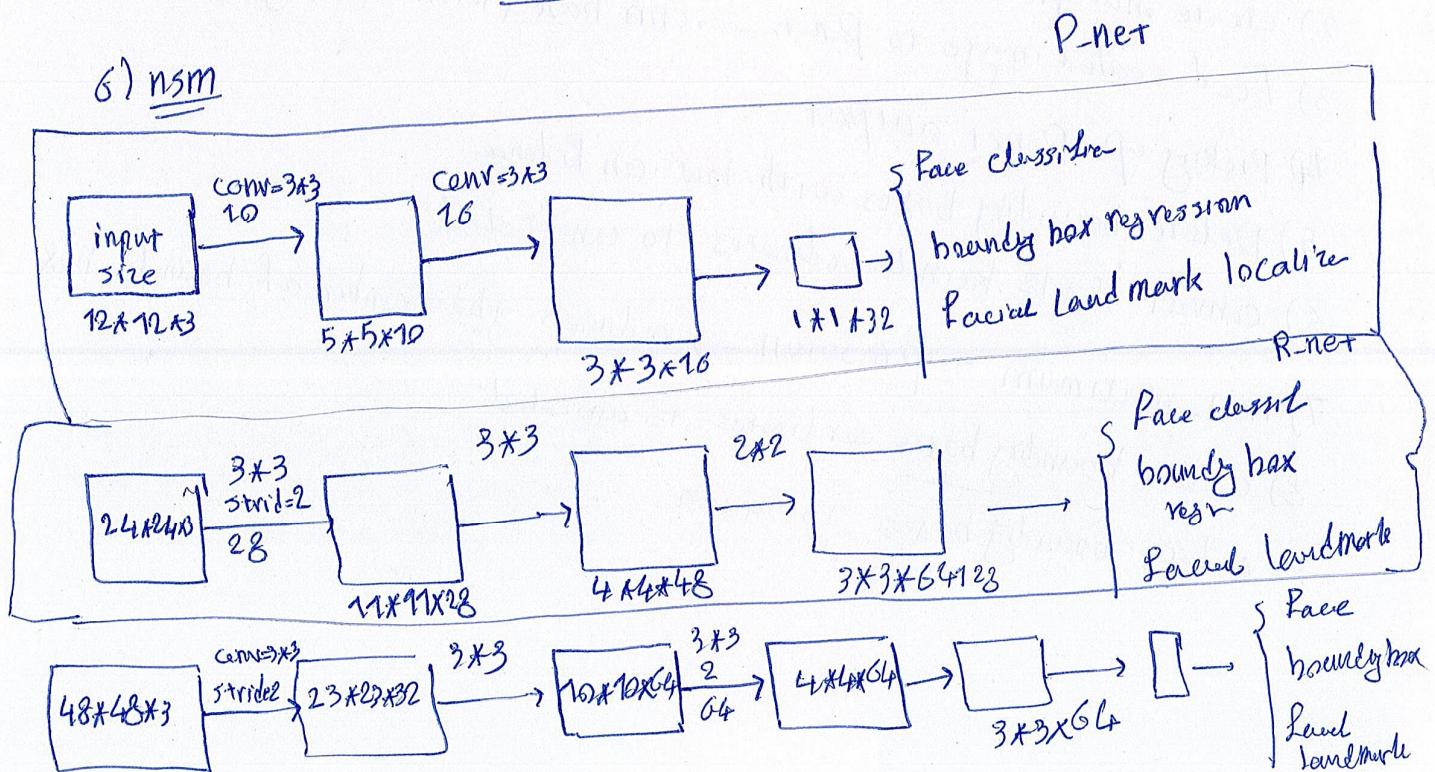
2) Feed scaled images into O-net

3) gather O-Net output

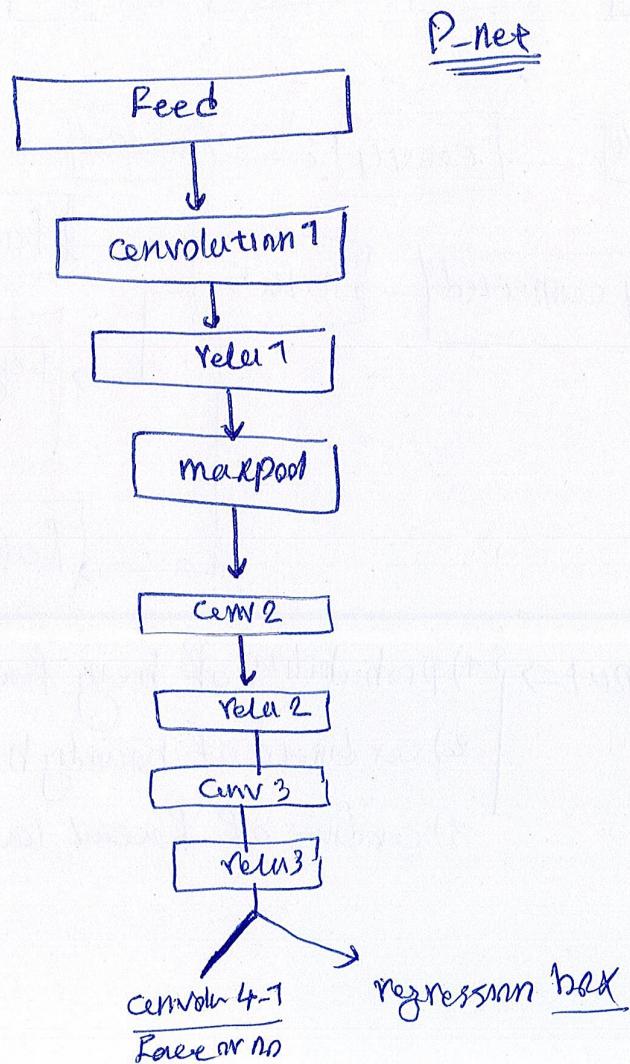
4) ~~concentration~~ Delete boundary boxes and boxes ~~and~~ with lower confidence

5) convert them to unscaled

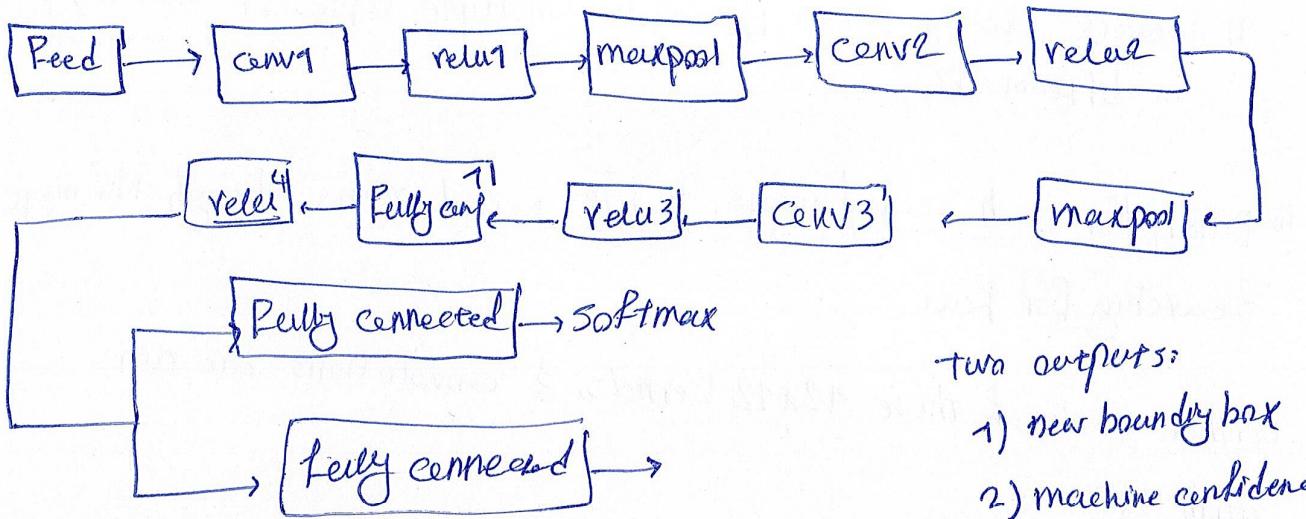
6) NMS



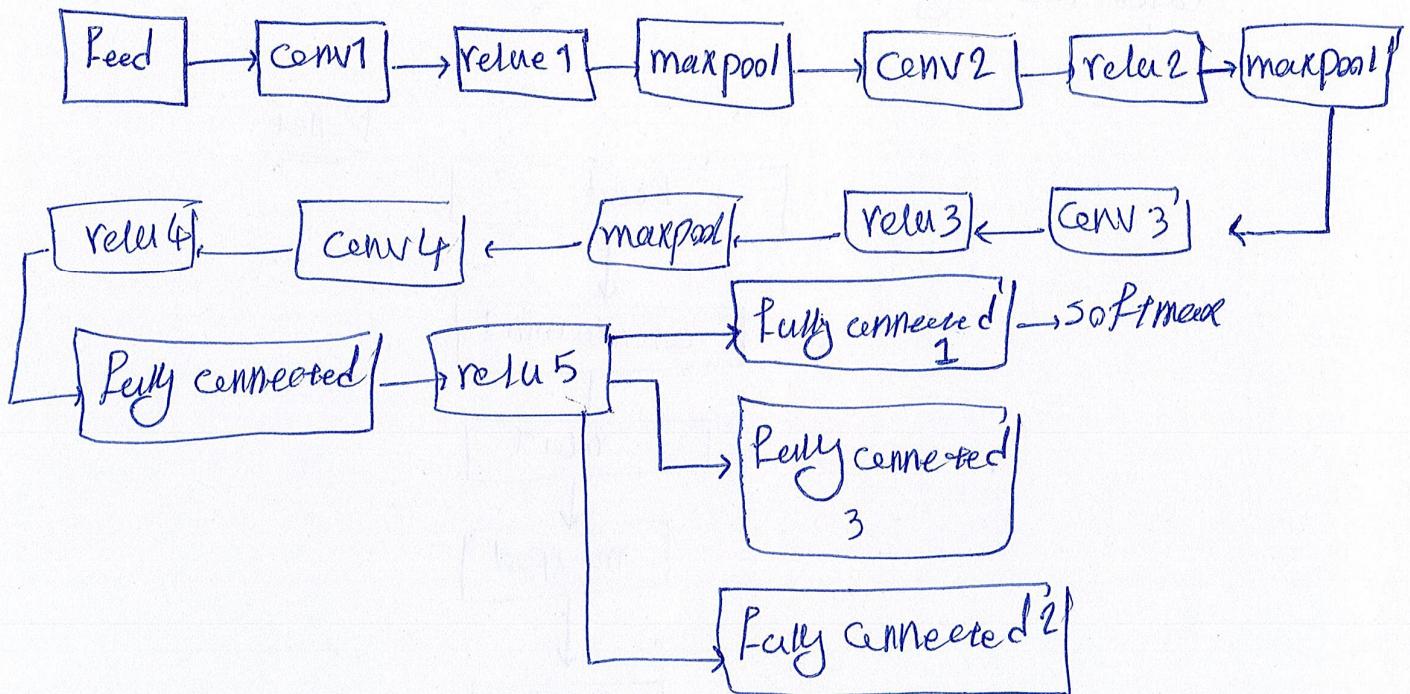
- For every image we pass in
  - 1) network creates Image Pyramid  $\rightarrow$  multiple copies of that image in different size.
- in P-net, For each scaled image,  $12 \times 12$  kernel runs through the image searching for Face.
  - within each of these  $12 \times 12$  kernels, 3 convolutions are run with  $3 \times 3$
  - after every convolution layer  $\rightarrow$  a relu layer  $\rightarrow$  as activation
  - then maxpooling  $\rightarrow$  reducing the size of input and take the largest.



## K-net



## O-net



3 output  $\Rightarrow$

- 1) probability of being face
- 2) coordinate of boundary box
- 3) coordinate of detected landmark

## Classification algorithms

↳ classical approach

↳ neural network base approach → nnb

- classical approaches → Do not use neural network → most of them around 99%. ~~and~~ were released before 2012.
- neural network base : using neural network as core → CNN

### ★★ classical approaches:

- For many years, they were the best option for classifying Images, ...
- with development of computer capability (cpu, Ram, GPU, ...) and new structure, they became weak → new methods took place for some reason

↳ 1) complexity / 2) great result

- SVM : a Famous, interesting and Flagship when it comes to talk about classification  $\Rightarrow$  40%

### Dynamic Bayesian network:

a summary information:

SVM: 40%, Dynamic Bayesian: 10%, Fuzzy: 6,7%.

marker model: 5%, KNN: 5, naive bias  $\Rightarrow$  5%, ...

Decision tree: 5%, adaboost: 5%, ...

## Some important dataset for facial expression

1) Affectnet  
more than one million images  
1250 emotion → tags  
Different language

2) Fer-2013  
28,000 labelled images → training set  
3,500      n      n → validation  
3,500      n      n → test set  
- collecting by Google Image  
- 7 emotion tag → { happy  
                  sad  
                  angry  
                  afraid  
                  surprise }

3) JAFFE  
- Japanese Female Facial Expression

- 1998
- around 10 Japanese women took picture of herself
- happiness, sadness, surprise, anger, disgust, neutral

4) MMI

- Create from 25 people of different regions. (Asia, European, etc.)  
- 44% woman, 55% men  
- has many non-basic feature → more than 6 labels.  
- each image is 720 × 576

## $\hookrightarrow$ SVM approaches

- 1) the most famous one - most used one, the most reliable one
- 2) SVM → is supervised learning.
- 3) capable of generating input-output mapping function.  
From a set of labeled data with features → produce function
- 4) classify element by determining boundaries (HyperPlane)
- 5) 1963 was developed  
separate class

## 2) Dynamic Bayesian network:

## 3) neural network base approaches:

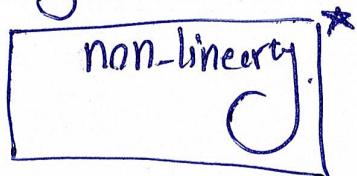
- 3.1) the use of CNN to solve computer vision has been successful in many areas.
  - 3.2) neural network base  $\xrightarrow{\text{MLP}}$   $\xrightarrow{\text{NN}}$   $\xrightarrow{\text{CNN}}$   
→ summary → for computer vision  
task → CNN: 66%  
neural netw 27,2  
mlp: 18,1%
  - 3.3) ANN (Artificial neural network):
    - they all attempt to simulate what happens in human brain.
    - they consist of processing unit → node
    - all nodes are interconnected with nodes in next layer → via connection → weight (synapse)
- \* mlp and CNN are  $\Rightarrow$  a specific configuration.

### 3.4) MLP: multi layer perceptron

- using some layer of node to improve degree of accuracy and extract complex pattern and features.

- hidden layer → great tools.

- sigmoid and relu function are used to perform



### 3.5) Convolutional Neural network: (C)

\* they are great to deal with Image/video

\*) take Image/video directly → Doesn't need to flat your image

\*) reduce time

\*) number of parameter decreased sharply

\*) CNNs have 2 important section

| - Feature extraction  
| - Fully connected

- Feature extraction: using convolution idea to extract feature

- Fully connected → MLP to classify Feature map.

## Feature selection:

- is one of most important stages for extracting important and relevant feature to have great prediction.
- there are several algorithms for Feature extraction in traditional machine learning
  - 1) Active Shape Model (ASM)
  - 2) Local Binary Patterns (LBP)
  - 3) HOG  $\Rightarrow$  Histogram of oriented gradient

1) ASM: - It was proposed in 1995

- with the goal of generating points that fit the ~~exact~~ contour of the object.

2) LBP, HOG:

- work by splitting the image into small pieces called cells and generating histogram.

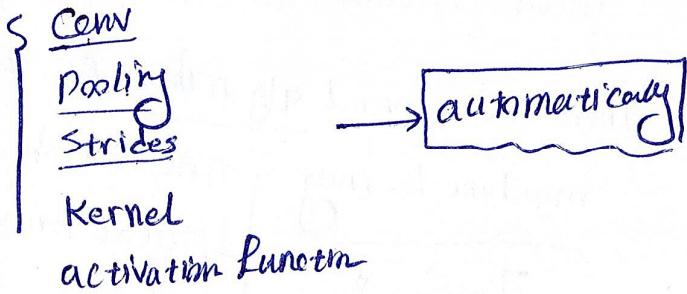
LBP: provides a comparison between each cell's pixels and their eight neighbors to build a binary number.

HOG: instead of making comparison between pixels, makes a histogram from the gradients.

## Feature extraction (CNN):

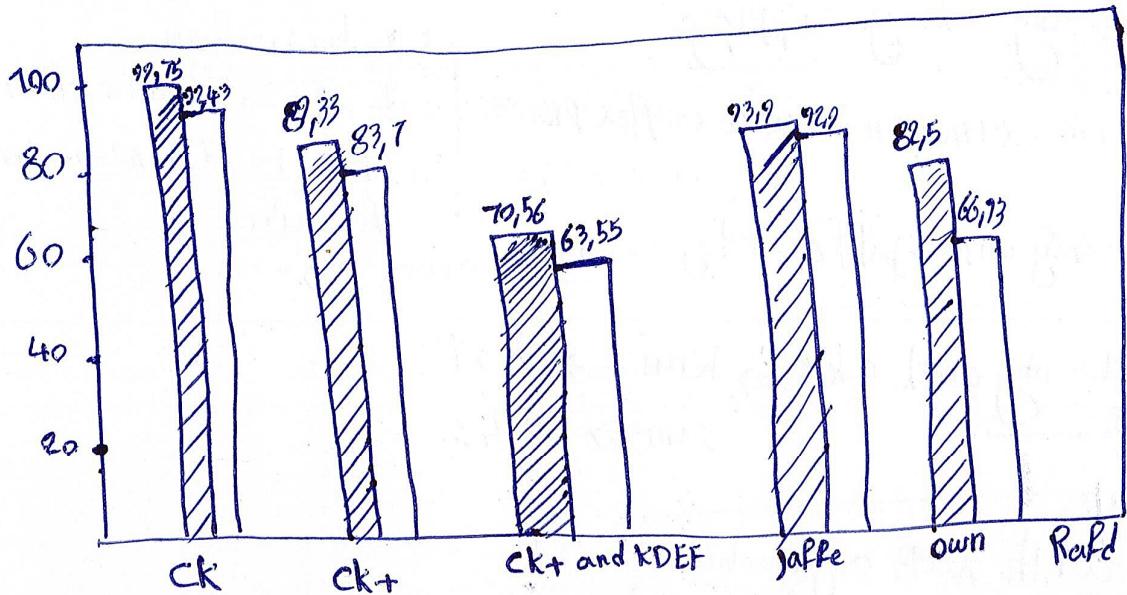
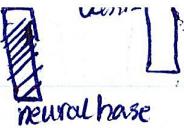
- CNN Does it in first section (Conv layer)

- Conv layer consists of



After passing Conv layer → output is sent to  
Fully connected to make prediction.

## emotion Recognition method.



\*) in order to achieve good result with CNN model, training dataset must be huge and large.

\*) data augmentation → great tool for generating Picture

\*) when dataset has less data → CNN model will overfitted. X

that is why some classical approach got great result. \*\*\*\*

## Jaffe with classical approach.

Preprocessing: Image cropping.

Feature extraction: quite complex process

training with Jaffe → they were able to achieve 100

testing with CK+ ⇒ Knn → 99.31

SVM ⇒ 99.43.

Quite complex XX

- 1) radon transform
- 2) EMD → minimize effect noise.
- 3) IMF: based on relation between features.

## Jaffe with NNB approaches:

CNN → achieved 100%.

→ there were three datasets → training → huge amount of data

jaffe → 2213	training
KDEF → 4900	
S FEW → 700	

Viola-Jones → detecting

\* huge amount of Images → stay away from overfitting → 100

testing → 100 → JAFFE

test → KDEF → 80

## CK+/BU-3DFE NNB approaches:

96,76 → CK

using data augmentation

91,89 → BU-3DFE